# A Constraint-Based Approach to Automatic Data Partitioning

Wonchan Lee
Computer Science Department
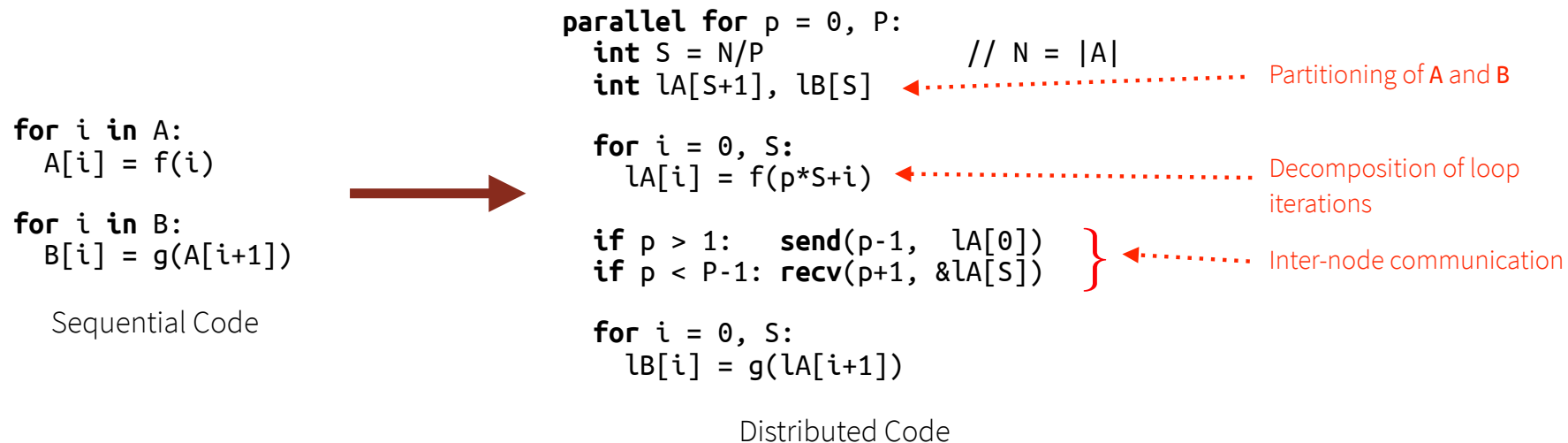Stanford University

# Auto-Parallelizers for Distributed Systems

- Goal: automatically generate distributed memory code from sequential programs

- Focus on data parallel programs

- Numerous efforts in the past several decades

  - High Performance Fortran and its predecessors (Fortran D, Vienna Fortran)

  - Polyhedral compilers for distributed memory machines

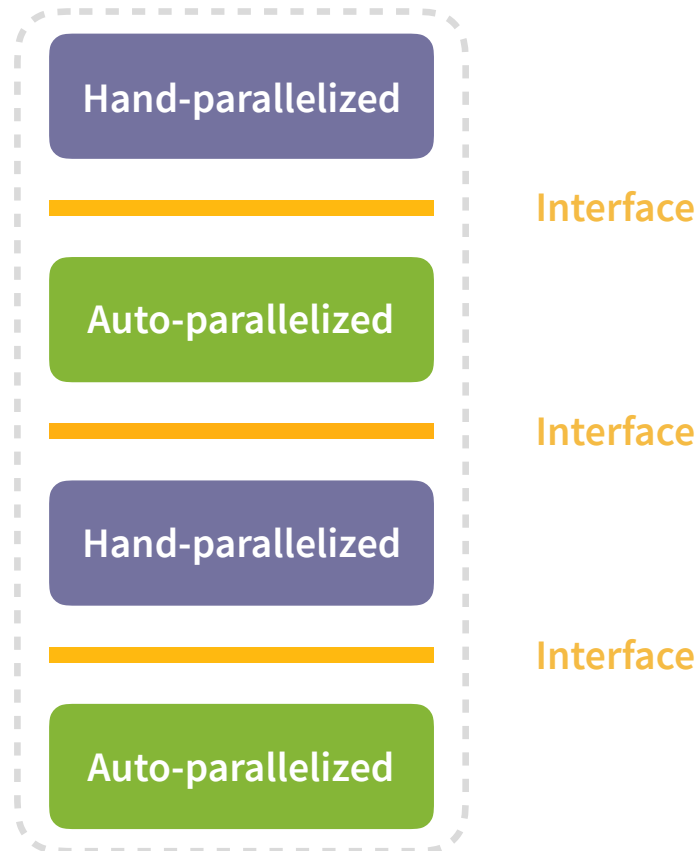**Why are they not being used more widely?**

2

# Issue 1: Configurability

- Compilers have to make decisions without the information only available at runtime

```
                              parallel for p = 0, P:
                                int S = N/P              // N = |A|
                                int lA[S+1], lB[S]   ⟵····· Partitioning of A and B

for i in A:                     for i = 0, S:
  A[i] = f(i)                     lA[i] = f(p*S+i)    ⟵····· Decomposition of loop
                                                              iterations
for i in B:                     if p > 1:   send(p-1,  lA[0])  }
  B[i] = g(A[i+1])              if p < P-1: recv(p+1, &lA[S])  }  ⟵····· Inter-node communication

     Sequential Code           for i = 0, S:
                                  lB[i] = g(lA[i+1])

                                        Distributed Code
```

- Decomposition of program and data must be determined at compile time, often by **hard-coded heuristics** in the compiler

- **Indirect accesses** make the output program even more obscure (inspect/executor)

# Issue 2: Composability

In practice, programs look like this:

**Hand-parallelized**

Interface

**Auto-parallelized**

Interface

**Hand-parallelized**

Interface

**Auto-parallelized**

We need **interface** for seamless integration,
but auto-parallelized parts are **opaque** to the rest of the program

# Data Distributions in HPF

- Annotation language to describe the primary partition of data

  - E.g., tiling on the first dimension of A:

    ```
              REAL A(1000,10000)
    !HPF$ DISTRIBUTE A(BLOCK,*)
    ```

- Can serve as an interface for both configuration and composition

  - Support for sharing data partitions is key to configurability and composability

- Limited because "data distributions were not themselves data objects"[†]

† Ken Kennedy, Charles Koelbel, and Hans Zima. 2007. The Rise and Fall of High Performance Fortran: An Historical Object Lesson. In Proceedings of the Third ACM SIGPLAN Conference on History of Programming Languages. ACM, 7–1.

# Programming Models with First-Class Partitions

Use **data partitions** as **programmable objects**

```
// pA1 is a partition of A
parallel for x in pA1:
  A = pA1[x]
  for i in A:
    A[i] = f(i)
```

Examples: **Legion**, StarPU, PaRSEC

**pA1** is an abstraction over partitions of **A** constructed at runtime:



**Data partitions** can naturally serve as an interface between different parts

# Programming Models with First-Class Partitions

Can provide synchronization and communication from multiple data partitions

```
// pA1 and pA2 are partitions of A
parallel for x in pA1:
  A = pA1[x]
  for i in A:
    A[i] = f(i)


parallel for x in pB:
  A = pA2[x]
  B = pB[x]
  for i in B:
    B[i] = g(A[i+1])
```

Implicit communication and synchronization for every i and j such that pA1[i] ∩ pA2[j] ≠ ∅

→ Can be handled automatically by compiler[‡] or runtime system[†]

† Michael Bauer, Sean Treichler, Elliott Slaughter, Alex Aiken, *Legion: expressing locality and independence with logical regions*. SC12.

‡ E. Slaughter, W. Lee, S. Treichler, W. Zhang, M. Bauer, G. Shipman, P. McCormick, and A. Aiken, *Control replication: Compiling implicit parallelism to efficient SPMD with logical regions*, SC17.

7

# Auto-Parallelization as Constraint Solving

Auto-parallelization amounts to finding **legal partitions** by solving **partitioning constraints**

```
parallel for x in pA1:
  A = pA1[x]
  for i in A:        ①
    A[i] = f(i)


parallel for x in pB:
  A = pA2[x]
  B = pB[x]
  for i in B:        ②
    B[i] = g(A[i+1])    ③
```

Find partitions **pA1**, **pA2**, and **pB** that satisfy these **constraints**:

① **pA1** covers **A**

② **pB** covers **B**

③ For any index **i** in **pB[x]**, **pA2[x]** includes **i+1**

# Constraint-Based Automatic Data Partitioning

Parallelizes sequential program using **data partitions**

Infers **partitioning constraints**

Discharges constraints with **interface constraints**

```
...
// Hand-parallelized code
...
assert(π(some_pA))
```

```
for i in A:
  A[i] = f(i)
```

```
require(π(pA))
parallel for x in pA:
  A = pA[x]
  for i in A:
    A[i] = f(i)
```

Or, **synthesizes partitioning code** using constraint solver

```
pA = some_pA
parallel for x in pA:
  ...
```

```
pA = partition(A,...)
parallel for x in pA:
  ...
```

# DPL as Constraint Language

- DPL(Dependent Partitioning Language)[†]: domain specific language for data partitioning

    - DPL programs construct data partitions using high-level operators

    - DPL operators have well-defined semantics and scalable implementation

- DPL can be used to describe both partitioning constraints and their solutions

† Sean Treichler, Michael Bauer, Rahul Sharma, Elliott Slaughter, and Alex Aiken, *Dependent Partitioning,* OOPSLA16.

Example:

```
parallel for x in pB:
  A = pA2[x]
  B = pB[x]
  for i in B:                ③
    B[i] = g(A[i+1])
```

③ $\forall j, \forall i \in pB[j], (i+1) \in pA2[j]$

A function that
maps **i** to **i+1**



Partition of the range of $\lambda i.i+1$

Collecting image of $\lambda i.i+1$

Partition of the domain of $\lambda i.i+1$

DPL program: **pA2 = image(pB, $\lambda i.i+1$)**

Example:

Many partitions can satisfy the constraint

```
parallel for x in pB:
    A = pA2[x]
    B = pB[x]
    for i in B:
        B[i] = g(A[i+1])
```

$\forall j, \forall i \in pB[j], i+1 \in pA2[j]$



Constraint that characterize **all legal partitions** for **pA2**:

$$\text{image}(pB, \lambda i.i+1) \subseteq pA2$$

The program **pA2 = image(pB, λi.i+1)** is one solution of this constraint

# Example: Particle Simulation

Updates the position of every particle using the velocity of cells

```
for i in PT:
  c = PT[i].cell
  PT[i].pos = f(C[c].vel,C[n(c)].vel)
```



PT[·].cell

n

PT (ParTicles)          C (Cells)

# Constraint Inference

Identifies necessary data partitions

Needs two partitions **pC1** and **pC2** of `C`
for different access patterns

```
for i in PT:
  c = PT[i].cell
  PT[i].pos = f(C[c].vel,C[n(c)].vel)
```

Needs a partition **pPT** of `PT`

**pPT[j]**                    **pC1[j]**

**pC2[j]**

$\forall j, \forall i \in \texttt{pPT[j]}, \texttt{PT[i].cell} \in \texttt{pC1[j]}$

$\forall j, \forall i \in \texttt{pC1[j]}, \texttt{n(i)} \in \texttt{pC2[j]}$

Infers partitioning constraints

Must run all iterations
at least once

```
require(complete(pPT,PT))
require(image(pPT,PT[·].cell) ⊆ pC1)
require(image(pC1,n) ⊆ pC2)
parallel for x in pPT:
  PT = pPT[x]; C1 = pC1[x]; C2 = pC2[x]
  for i in PT:
    c = PT[i].cell
    PT[i].pos = f(C1[c].vel,C2[n(c)].vel)
```

A function that maps
particles to cells

# Solving Constraints

Two solutions for **require(image($\text{pPT}$,PT[·].cell) ⊆ pC1)** :



✔ **pPT** = **equal**(PT,N)
 **pC1** = **image**(**pPT**,PT[·].cell)

**pC1** = **equal**(C,N)
 **pPT** = **preimage**(PT[·].cell,**pC1**)

**require**(**complete**(**pPT**,PT))
**require**(**image**(**pPT**,PT[·].cell) ⊆ **pC1**)
**require**(**image**(**pC1**,n) ⊆ **pC2**)

Solve
constraints
→

**pPT** = **equal**(PT,N)
**pC1** = **image**(**pPT**,PT[·].cell)
**pC2** = **image**(**pC1**,n)

# Handling Multiple Loops

One loop = One set of partitioning constraints

```
for i in PT:
  c = PT[i].cell
  PT[i].pos = g(C[c].vel,C[n(c)].vel)
```

```
require(complete(pPT,PT))
require(image(pPT,PT[·].cell) ⊆ pC1)
require(image(pC1,n) ⊆ pC2)
```

```
for i in C:
  C[i].vel  = h(C[i].acc,C[n(i)].acc)
```

```
require(complete(pC3,C))
require(image(pC3,n) ⊆ pC4)
```

✔ Capture **all possible partitioning strategies**

✘ Can lead to **excessive communication** if solved naïvely

# Handling Multiple Loops

Constraint solver **unifies** partitions to maximize partition reuse

```
for i in PT:
  c = PT[i].cell
  PT[i].pos = g(C[c].vel,C[n(c)].vel)
```

```
require(complete(pPT,PT))
require(image(pPT,PT[·].cell) ⊆ pC1)
require(image(pC1,n) ⊆ pC2)
```

```
for i in C:
  C[i].vel  = h(C[i].acc,C[n(i)].acc)
```

```
require(complete(pC3,C))
require(image(pC3,n) ⊆ pC4)
```

<center>Similar access patterns</center>

<center>Isomorphic constraints</center>

<center>Unified!</center>

```
          pC3 = pC1 = equal(PT,N)
Solution: pC4 = pC2 = image(pC1,n)
          pPT = preimage(PT[·].cell,pC1)
```

# External Constraints

In the real simulation code, particles might **move to different cells**,

```
while t < T:
  for i in PT:
    c = PT[i].cell
    PT[i].pos = g(C[c].vel,
                  C[n(c)].vel)

  for i in C:
    C[i].vel = g(C[i].acc,
                 C[n(i)].acc)

  for i in PT:
    PT[i].cell = h(PT[i].pos)
```

Solve constraints →

requiring **pPT** to be **repartitioned** every time step ☹

```
pC1 = equal(C,N)
pC2 = image(pC1,n)

while t < T:
  pPT = preimage(PT[·].cell,pC1)

  parallel for x in pPT:
    ...

  parallel for x in pC1:
    ...

  parallel for x in pPT:
    ...
```

18

# External Constraints

User can **manually parallelize** particle transfer code
and provide external constraints as an interface:

```
pC3 = pC1 = pCell
pC4 = pC2 = image(pCell,n)
```

```
while t < T:
  ...
  // Manual particle transfer code using
  // pParticle and pCell
  ...
  assert(
    image(pParticle,PT[·].cell) ⊆ pCell))

  ...
```

Solve constraints with
external constraints →

```
while t < T:
  ...
  // Manual particle transfer code
  ...
  pPT = pParticle
  ...
```

Unifiable constraints

No more repartitioning 🙂

Partitioning constraints:

```
require(complete(pPT,PT))
require(image(pPT,PT[·].cell) ⊆ pC1)
require(image(pC1,n) ⊆ pC2)
require(complete(pC3,C))
require(image(pC3,n) ⊆ pC4)
```

External constraints provide a precise control
over the automated data partitioning process

# Evaluation

- Implemented the constraint inference and solver algorithms in Regent[†]

  - Regent is a high-level programming language with **first-class data partitions and DPL**

- Weak scaling performance of four benchmark programs

  - Stencil: 9-point stencil in 2D grid

  - MiniAero: explicit Navier-stokes solver on hexahedral 3D mesh

  - Circuit: circuit simulator on unstructured circuit graphs

  - PENNANT: Lagrangian hydrodynamics on unstructured 2D mesh

- **Machine: Piz Daint** (12-core Xeon E5-2690, NVIDIA P100, and 64 GB memory per node)

- All benchmark programs ran on GPUs

† Elliott Slaughter, Wonchan Lee, Sean Treichler, Michael Bauer, Alex Aiken, *Regent: a high-productivity programming language for HPC with logical regions.* SC15.

# Weak Scaling: Stencil and MiniAero

Stencil ($0.9 \times 10^9$ points/node, Piz Daint)

MiniAero ($2.1 \times 10^6$ Cells/node, Piz Daint)



Auto-parallelized programs match hand-parallelized programs within 3%

Circuit ($100 \times 10^3$ wires/node, Piz Daint)

Circuit ($100 \times 10^3$ wires/node, Piz Daint)

Add interface constraints

Auto-Parallelized Circuit

**Parallel Circuit Generator**

**Auto-Parallelized Compute Tasks**
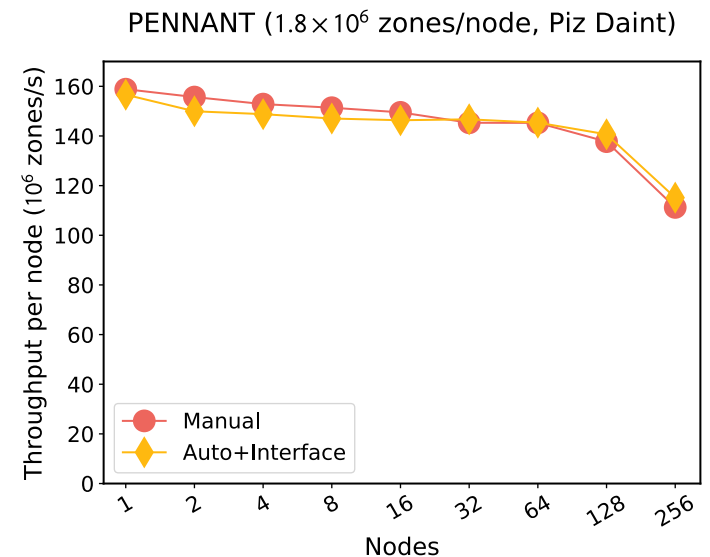
Generator uses two node partitions:
pNodes_private and pNodes_shared

Interface constraints:
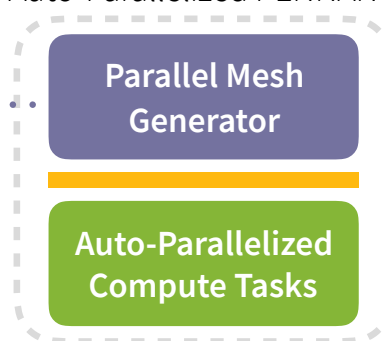complete(pNodes_private ∪ pNodes_shared, Nodes)

22

PENNANT ($1.8 \times 10^6$ zones/node, Piz Daint)



**Add interface constraints**

PENNANT ($1.8 \times 10^6$ zones/node, Piz Daint)



Two points partitions:
pPoints_private and pPoints_shared



Each side of a polygon colocates with
its previous and next sides

Auto-Parallelized PENNANT

**Parallel Mesh Generator**

**Auto-Parallelized Compute Tasks**

**Interface constraints:**

complete(pPoints_private ∪ pPoints_shared, Points)
image(pSides, Sides[·].prev_side) ⊆ pSides
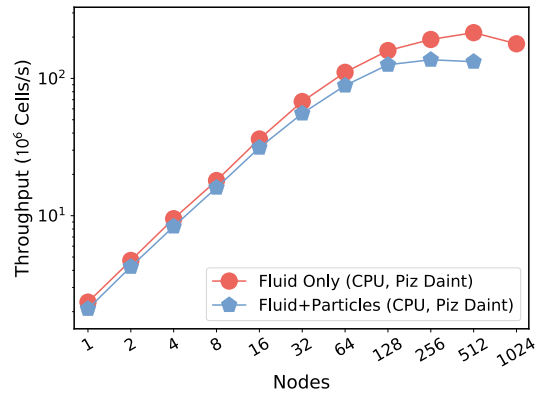image(pSides, Sides[·].next_side) ⊆ pSides

23

# Case Study: Soleil-X

- Developed for the PSAAP II program at Stanford

- Eulerian Fluid + Lagrangian Particles + Radiation (DOM/Algebraic)

  - DOM is manually parallelized

  - Fluid and particles are auto-parallelized except for particle transfers

Heated section of concentrated
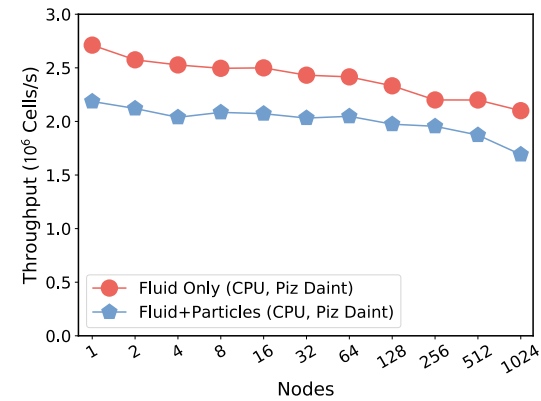solar energy receiver

# Soleil-X Performance



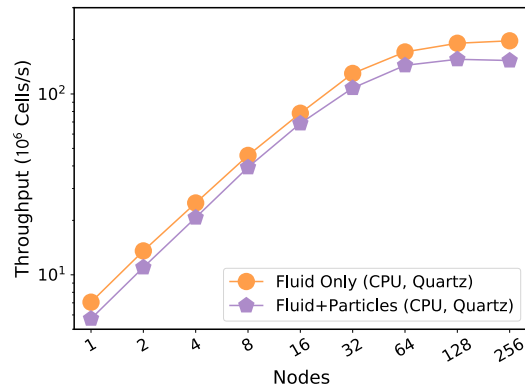Strong Scaling ($512^2 \times 256$ Cells, 4M Particles)

92X speedup at 512 for fluid
65X speedup at 256 with particles

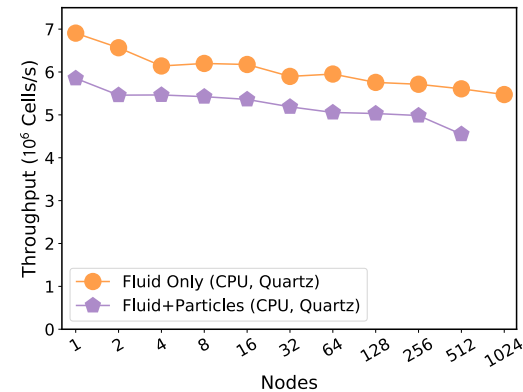Weak Scaling ($256^3$ Cells/Node, 1M Particles/Node)

77% parallel efficiency
(9,216 CPUs)

Strong Scaling ($512^3$ Cells, 8M Particles)

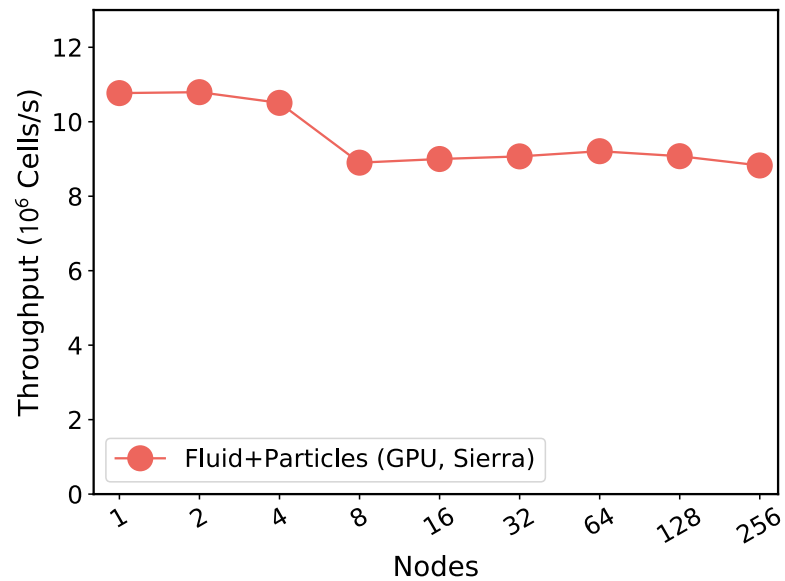17X speedup at 128 for fluid
15X speedup at 128 with particles

Weak Scaling ($512 \times 256^2$ Cells/Node, 2M Particles/Node)

79% parallel efficiency
(30,720 CPUs)

# Soleil-X Performance

Weak Scaling (67M Cells/Node, 32M Particles/Node)



82% parallel efficiency
(1,024 GPUs)

- Solves multi-component Navier-Stokes equations in compressible formulation

  - Accounts for complex chemistry and multicomponent transport

- Heavy flux tasks are auto-parallelized



**Hypersonic Task-based Research (HTR) solver**

**Transport equations**

$$\frac{\partial \rho_i}{\partial t} + \nabla \cdot (\rho_i \mathbf{u} + \rho_i \mathbf{V}_i) = \dot{\omega}_i, \quad \text{for } i = 1,\ldots,N_s$$

$$\frac{\partial (\rho \mathbf{u})}{\partial t} + \nabla \cdot (\rho \mathbf{u}\mathbf{u}) = \nabla \cdot \bar{\bar{\tau}} - \nabla p$$
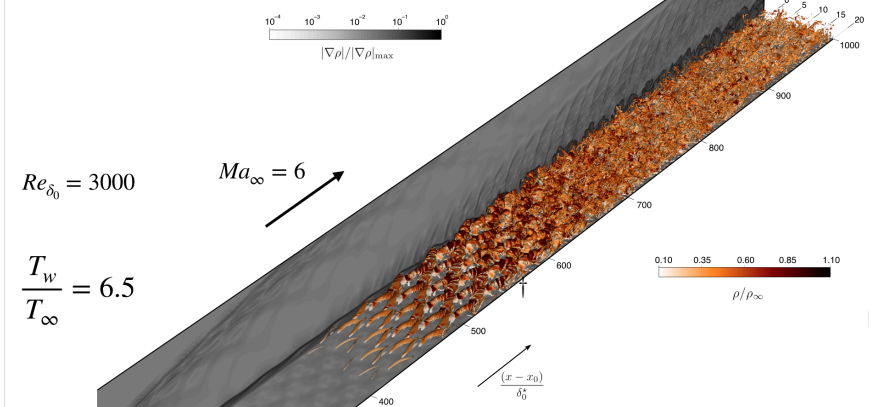
$$\frac{\partial (\rho E)}{\partial t} + \nabla \cdot (\rho \mathbf{u} H) = \nabla \cdot (\lambda \nabla T) + \nabla \cdot (\bar{\bar{\tau}}\mathbf{u}) - \sum_i^{N_s} \nabla \cdot (\rho_i \mathbf{V}_i h_i)$$

$$\mathbf{V}_i = -D_i \nabla \ln X_i + \sum_i^{N_s} Y_i D_i \nabla \ln X_i$$

$$\sum_i^{N_s} \rho_i = \rho = \sum_i^{N_s} \frac{pX_iW_i}{\mathcal{R}^0 T}$$

- TENO6 low-dissipation scheme for Euler fluxes
- Second-order scheme for diffusion fluxes
- Shock capturing capabilities
- Multicomponent transport
- Arrhenius chemistry computed at runtime with or without time-operator splitting
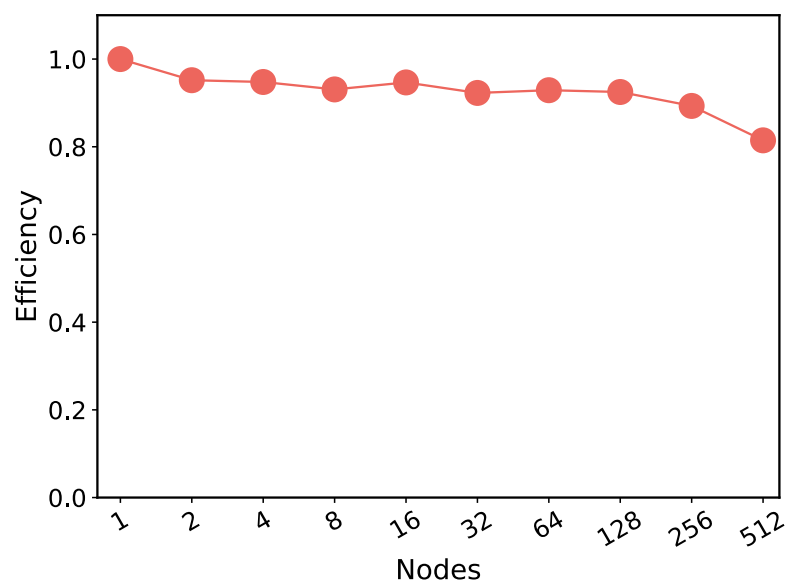- Thermo-physical and transport properties computed at runtime



**Hypersonic transitional boundary layer (low-enthalpy)**
**Reference: Franko and Lele (2013)**

$Re_{\delta_0} = 3000$

$Ma_\infty = 6$
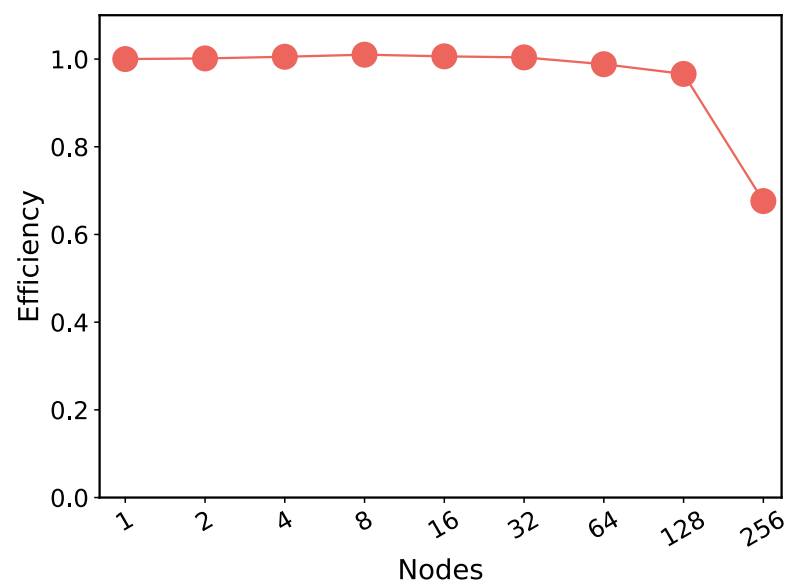
$\dfrac{T_w}{T_\infty} = 6.5$

# HTR Performance

Quartz (CPUs)



Lassen (GPUs)



81% efficiency
18,432 CPUs
400M points
1s/timestep

68% efficiency
1,024 GPUs
4.8B points
0.7s/timestep

28

# Conclusion

- First-class data partitions enable composable and configurable auto-parallelization

- A constraint-based data partitioning brings scalability of manual parallelization to auto-parallelized programs

# Questions?

**Legion**

*A Data-Centric Parallel Programming System*

Github

Legion is a data-centric parallel programming system for writing portable high performance programs targeted at distributed heterogeneous architectures. Legion presents abstractions which allow programmers to describe properties of program data (e.g. independence, locality). By making the Legion programming system aware of the structure of program data, it can automate many of the tedious tasks programmers currently face, including correctly extracting task- and data-level parallelism and moving data around complex memory hierarchies. A novel mapping interface provides explicit programmer controlled placement of data in the memory hierarchy and assignment of tasks to processors in a way that is orthogonal to correctness, thereby enabling easy porting and tuning of Legion applications to new architectures.

To learn more about Legion you can:

- Read the overview
- Visit the getting started page
- Download our publications
- Ask questions on our mailing list

**About Legion**

Legion is developed as an open source project, with major contributions from LANL, NVIDIA Research, SLAC, and Stanford. This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of two U.S. Department of Energy organizations (Office of Science and the National Nuclear Security Administration) responsible for the planning and preparation of a capable exascale ecosystem, including software, applications, hardware, advanced system engineering, and early testbed platforms, in support of the nation's exascale computing imperative. Additional support has been provided to LANL and SLAC via the Department of Energy Office of Advanced Scientific Computing Research and to NVIDIA, LANL and Stanford from the U.S. Department of Energy National Nuclear Security Administration Advanced Simulation and Computing Program. Previous support for Legion has included the U.S. Department of Energy's ExaCT Combustion Co-Design Center and the Scientific Data Management, Analysis and Visualization (SDMAV) program, DARPA, the Army High Performance Computing Research Center, and NVIDIA, and grants from OLCF, NERSC, and the Swiss National Supercomputing Centre (CSCS).

**Legion Contributors**

| Stanford | SLAC | LANL | NVIDIA |
|---|---|---|---|
| Todd Warszawski | Elliott Slaughter | Pat McCormick | Michael Bauer |
| Wonchan Lee | Alan Heirich | Samuel Gutierrez | (NVIDIA site) |
| Zhihao Jia | Seema Mirchandaney | Galen Shipman | Sean Treichler |
| Karthik Srinivasa Murthy | Seshu Yamajala | Jonathan Graham | |
| Manolis Papadakis | | Irina Demeshko | |
| Alex Aiken | | Nick Moss | |
| | | Wei Wu | |